

# Terascale Data Management

## Technology

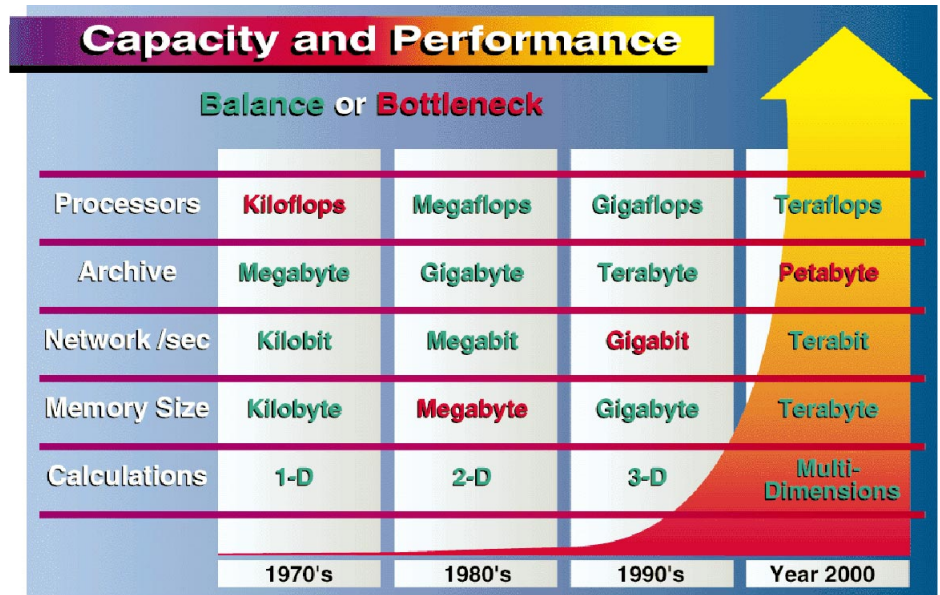
Efficient access to the massive amounts of data generated in a terascale computing environment requires a balanced, scalable architecture. Three major technologies supporting the Lawrence Livermore National Laboratory terascale data management environment include:

- Visualization.
- Scientific Data Management.
- High-performance storage.

## Problem Definition

The capability of computing technology has steadily increased for more than four decades, resulting in machines with very large memories, greater calculational performance, large disk caches, and higher-performing network interfaces. This increased capability, as shown in the chart on this page, has enabled applications to scale accordingly in calculational complexity, overall throughput, and resultant data.

Current data analysis techniques break down when operating in this increased capacity computing environment because of the massive amounts of data and the complexity of the data management problems. The data simply overwhelm the storage systems and networks. Therefore, terascale computing requires a new model for the way that we store, retrieve, and present data. The challenge is to provide a terascale data management environment that enables scientists to concentrate on science while minimizing the task of physically



Increased capacity enables applications to scale accordingly in complexity, throughput, and resultant data.

managing data and computer resources. Meeting this challenge requires the integration of several hardware and software components, including

- Data and information management,
- Visualization and analysis tools,
- High-performance storage and I/O, and
- High-performance networks.

## Technologies in the End-to-End Solution

Providing easy and fast access to the immense amount of data being generated requires a balanced, scalable architecture to help manage the data in the parallel and distributed computing environments. The architecture must address all of the levels in the end-to-end infrastructure, from the user interface through all of the hardware and software layers that support the data access and storage as shown in the infrastructure figure on the back page.

The LLNL strategy is to create an infrastructure that is flexible, extensible, and scalable so that the latest hardware and software can be integrated as technology advances. Commercial off-the-shelf technology is used wherever possible, and development projects leverage collaborations with vendors and universities.

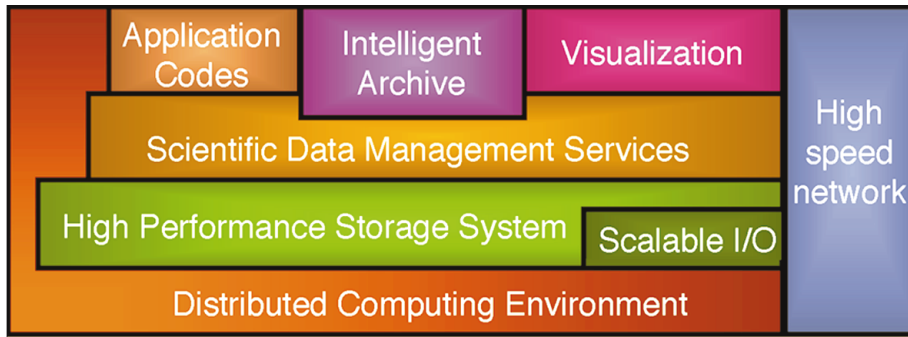
For example, terascale data management relies on technologies including high-speed networks and the Open Software Foundation's Distributed Computing Environment (DCE). Several development projects are under way at LLNL to address the areas of technology for which solutions are not commercially available.

## Scientific Data Management Services

Scientific Data Management (SDM) encompasses a number of technologies that can be integrated to accelerate, enhance, and simplify access to scientific data. An effective SDM environment would allow scientists to retrieve, store, and select data in the natural context of their work. It would seamlessly integrate databases, mass storage systems, networking, and computing resources to provide intelligent assistance in managing the complexity and scale of terascale data. The LLNL SDM project's strategy for developing an SDM environment is based on three levels of service:

- A large-capacity data storage manager.
- A data model supported by self-describing portable binary file formats or advanced database systems.

# Terascale Data Management



*The Terascale Data Management Infrastructure: End users and applications rely on terascale data management services to provide access to a wide range of computers and peripherals.*

- An interface to a suite of application tools to aid in data management and discovery.

## User Interface: Intelligent Archive

The Intelligent Archive (IA) project leverages metadata to help scientists organize, search, and interact with information and data. As the top level data management layer for end users, the IA provides graphical user interfaces and Web-based services to support quick and easy access to data. With the IA, scientists can not only access archived data but can also manage their current work and organize it for future access. Key strategies in the IA project include using commercial database software to manage metadata, writing tools in Java to ensure platform independence, coupling scientific code sys-

tems with IA tools, and supporting both automatic and interactive generation of metadata.

## Visualization: Distributed Desktop Delivery

The visualization component of terascale data management will provide techniques for effectively exploring huge multi-dimensional data sets. The Graphics Server project will deliver scientific visualizations to the desktop quickly and efficiently with the integration of high-performance graphics servers, storage systems, and massively parallel processors (MPPs) connected by a high-speed network. The range of services will include interactive visualization packages on individual desktops, memory-rich specialized graphics servers shared by several users, and renderers on massively

parallel platforms shared by large user communities.

## Hierarchical Storage Management: HPSS

The High Performance Storage System (HPSS) is a collaborative development project with IBM Government Systems, DOE laboratories, universities, vendors, and other research centers. The goal of the HPSS project is to improve the performance and capacity of hierarchical storage systems along with the architecture and functionality. HPSS provides a scalable, parallel storage system for highly parallel computers as well as traditional supercomputers and workstation clusters. HPSS requirements are driven by the high end of storage system and data management requirements and, although developed to scale for order-of-magnitude improvements, HPSS is a general-purpose storage system.

The Scalable I/O Facility (SIOF) project is developing software to enable I/O performance to scale with the computing performance and eliminate the current I/O bottleneck. The SIOF project is working with commodity disk and tape vendors to develop and supply networked attached devices that will be supported by the HPSS system.

## Related Information and Projects

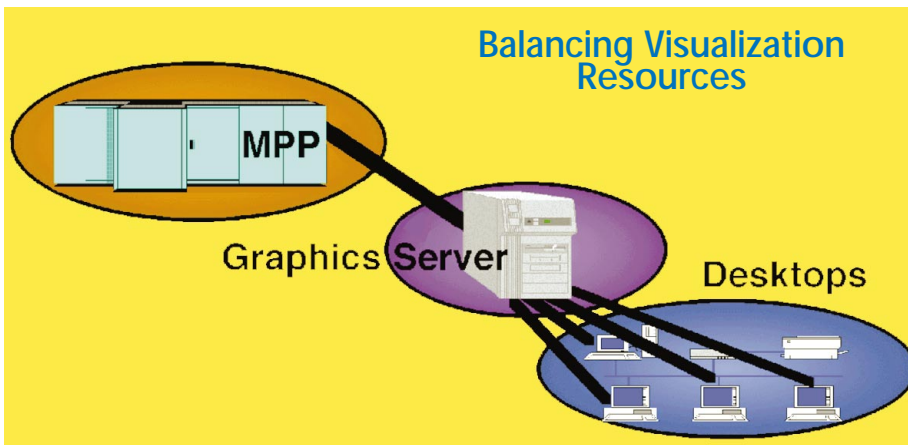
Livermore Computing  
[http://www.llnl.gov/liv\\_comp/lc](http://www.llnl.gov/liv_comp/lc)

Intelligent Archive Project  
[http://www.llnl.gov/liv\\_comp/ia](http://www.llnl.gov/liv_comp/ia)

HPSS  
<http://www.sdsc.edu/hpss/>

Scalable I/O Facility  
[http://www.llnl.gov/liv\\_comp/siof](http://www.llnl.gov/liv_comp/siof)

*For further information, contact Celeste Matarazzo, 510-423-9838, [celestem@llnl.gov](mailto:celestem@llnl.gov); Robyne Sumpter, 510-423-5054, [sumpter1@llnl.gov](mailto:sumpter1@llnl.gov); or Rebecca Springmeyer, 510-423-0794, [springme@llnl.gov](mailto:springme@llnl.gov).*



*The Distributed Desktop Delivery project will integrate high-performance graphics servers into the computing environment to perform visualization tasks that would otherwise fall to either the MPPs generating the data or the small systems on the desktops of the end-user scientists.*